**BMS**

# Next generation sequencing

생물정보팀 전형린

# Contents

- **1주차 Next generation sequencing 및 bioinformatics 개요**
  - High throughtput DNA decode
  - 세대 별 시퀀싱 차이
  - Next generation sequencing 원리
  - Application based on NGS
- **2주차 Hybrid-Seq 개요 및 실습**
  - Hybrid-Seq 원리
  - Hybrid-Seq 분석 과정 및 결과 설명
  - 분석 실습
- **3주차 R을 이용한 plotting**
  - ggplot2 을 이용한 plotting
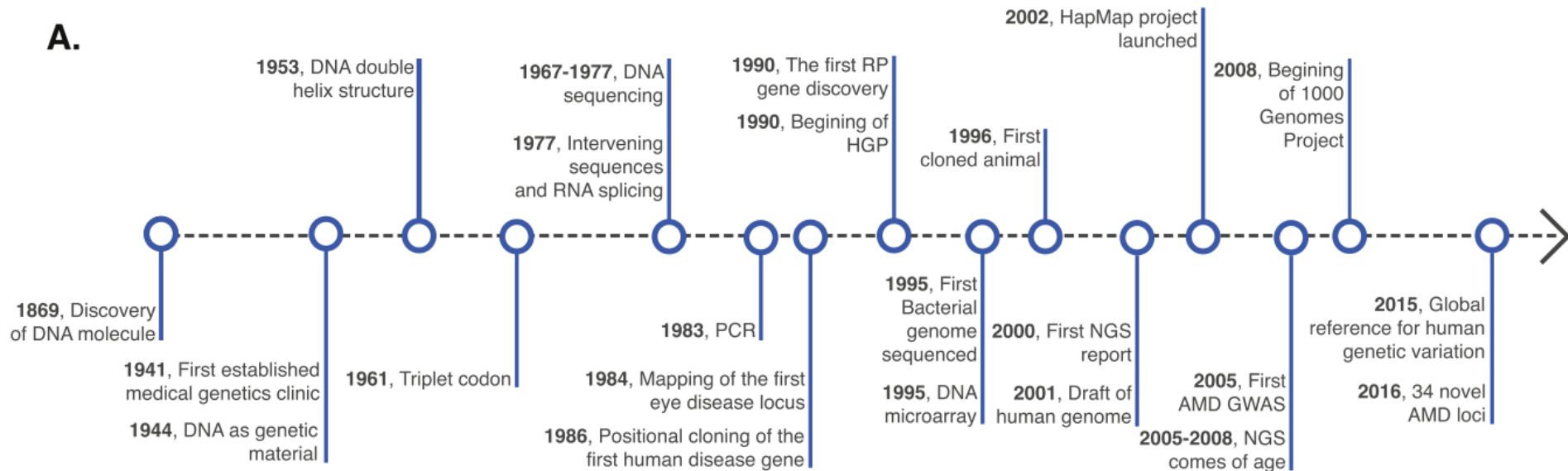  - R을 이용한 phylogenetic tree plotting

# Human Genome Project

# Timeline for Human Genome Project

# Timeline until HGP

- **Human genome project formally launched in 1990 and was declared complete on April 14, 2003**



*Nature*, 2001

# White House on 26 June 2000

BMS  Bio-Medical Science Co., Ltd.

# Cost for Human Genome Project

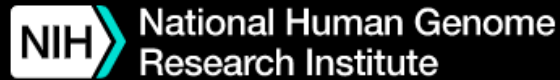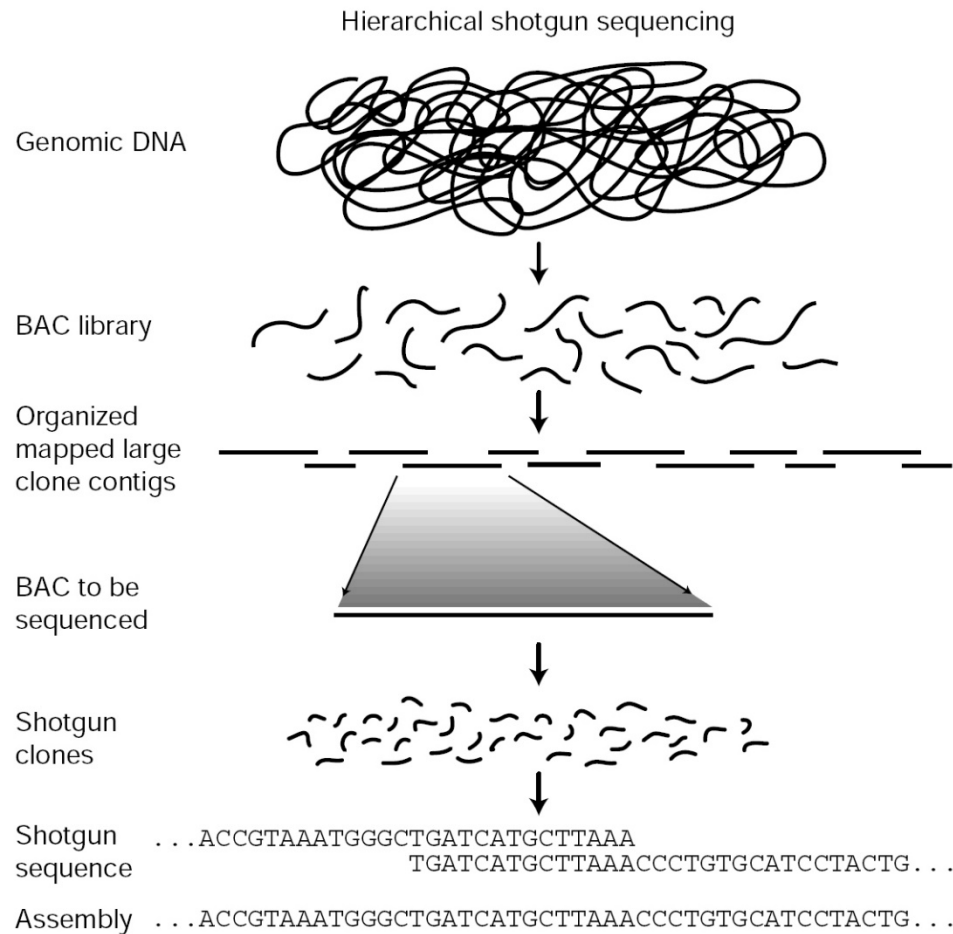### National Human Genome Research Institute

## How much did it cost?

In 1990, Congress established funding for the Human Genome Project and set a target completion date of 2005. Although estimates suggested that the project would cost a total of $3 billion over this period, the project ended up costing less than expected, about $2.7 billion in FY 1991 dollars. Additionally, the project was completed more than two years ahead of schedule.
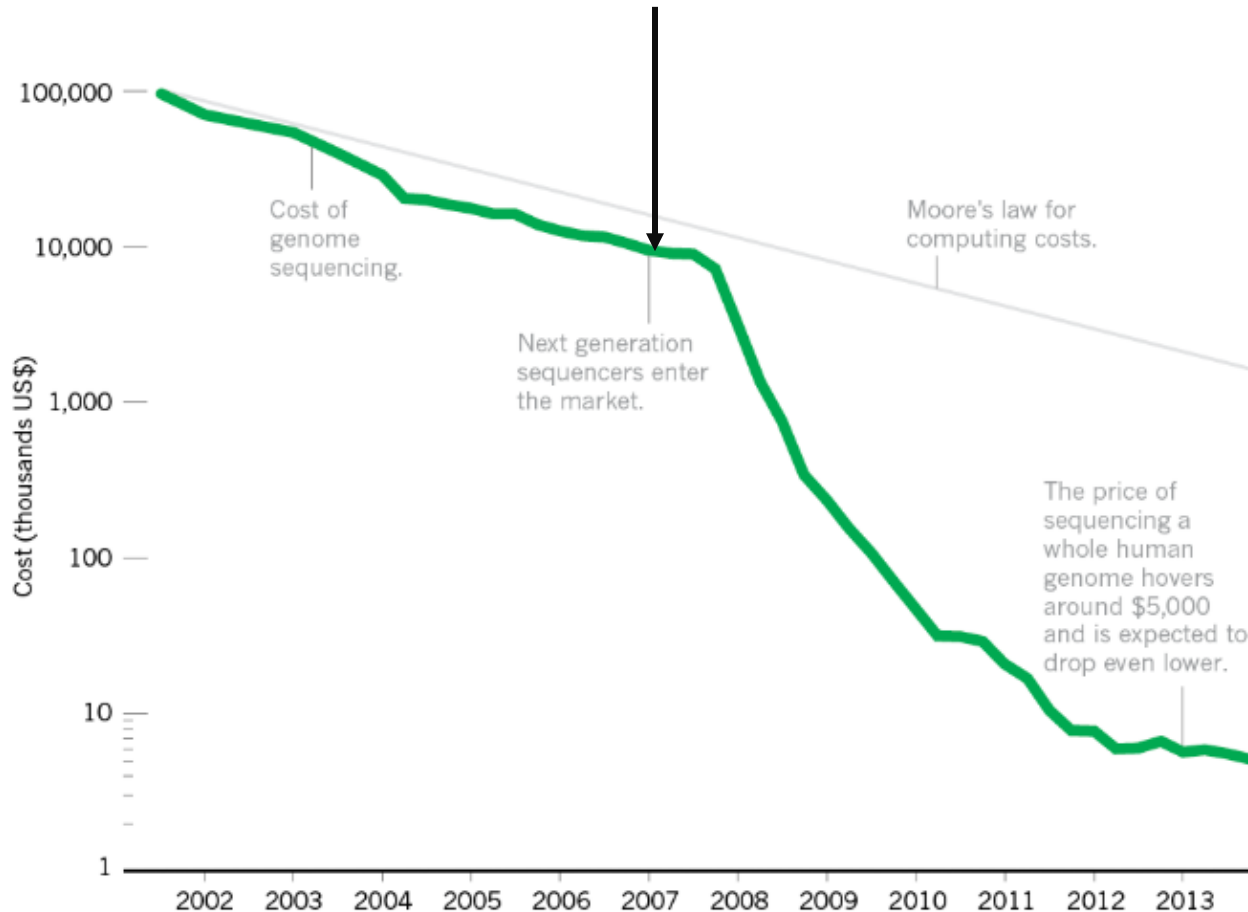
# Strategy for sequencing



Hierarchical shotgun sequencing

Genomic DNA

BAC library

Organized
mapped large
clone contigs

BAC to be
sequenced

Shotgun
clones

Shotgun
sequence    ...ACCGTAAATGGGCTGATCATGCTTAAA
                        TGATCATGCTTAAACCCTGTGCATCCTACTG...

Assembly   ...ACCGTAAATGGGCTGATCATGCTTAAACCCTGTGCATCCTACTG...

BMS  Bio-Medical Science Co., Ltd.

# Cost for genome sequencing

**Next Generation Sequencing platform was appeared as game exchanger in sequencing field**



*https://www.researchgate.net/figure/a-Cost-per-base-of-the-different-sequencing-techniques-as-a-function-of-time-The-gray_fig2_271772842*

# Next Generation Sequencing Platform

- Roche/454 (GS FLX+/GS Junior)
- Illumina Genome Analyzer (HiSeq/MiSeq/NextSeq)
- Life Technologies (3500 Genetic Analyzer)
- Ion Torrent Proton/PGM)
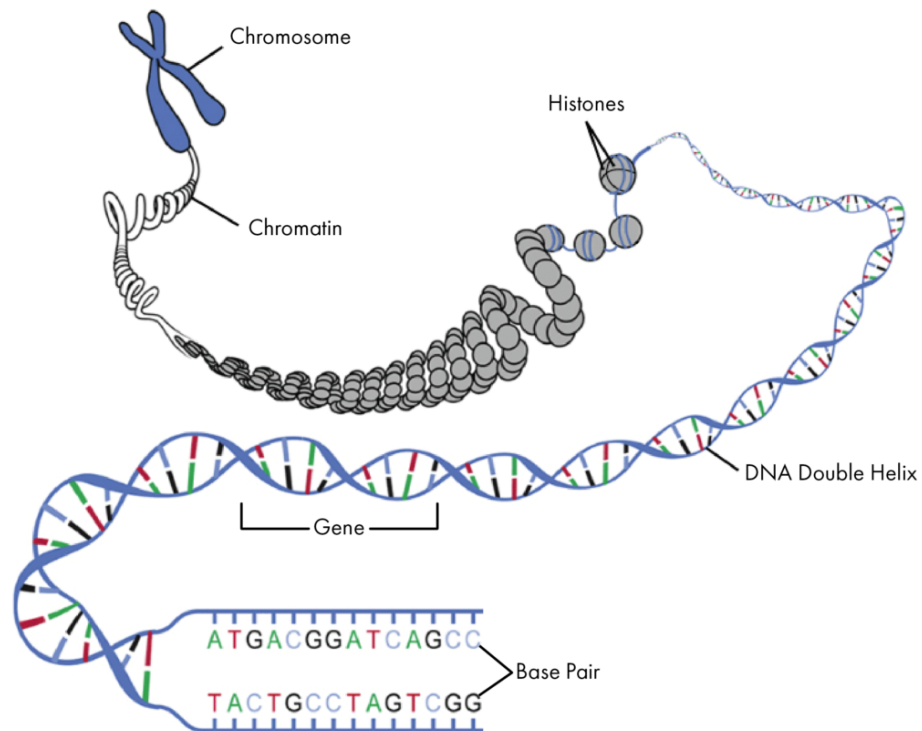- Applied Biosystems (SOLiD, 3730xl DNA Analyzer )

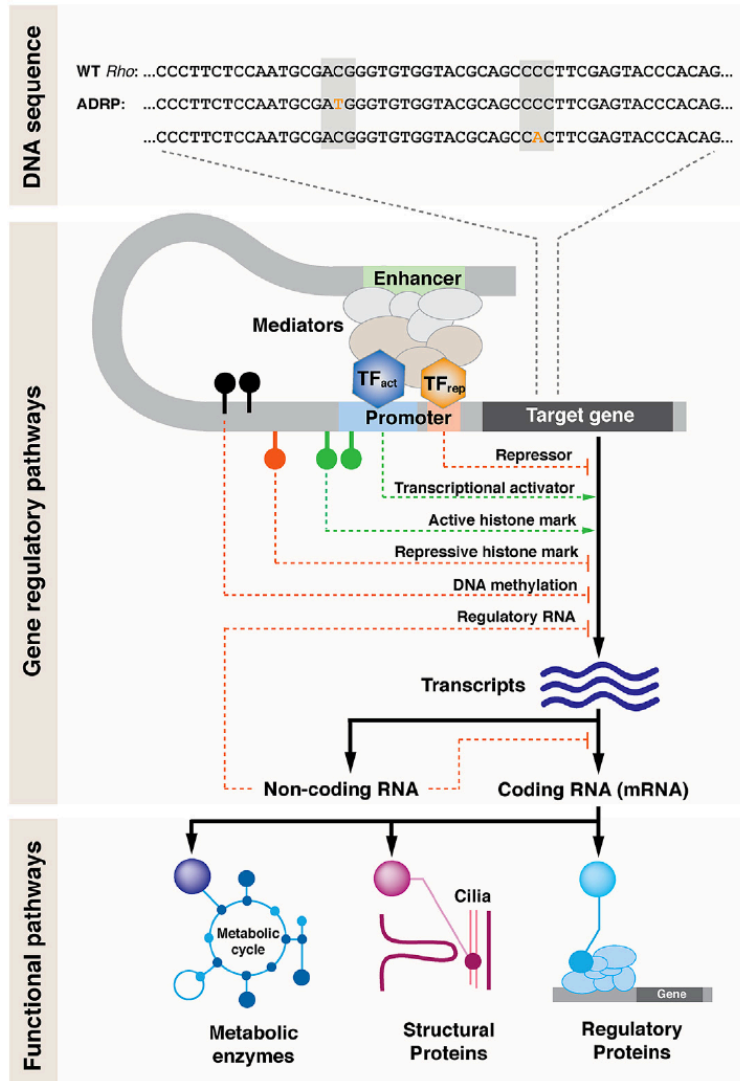*http://www.vib.be/en/about-vib/annual-report/2012/research/activities/Pages/Service%20Facilities.aspx*

# "DNA Sequencing"

- DNA sequencing involves the use of various methods for determining the order of the nucleotide bases — **adenine**, **cytosine**, **guanine**, and **thymine** — in a molecule of DNA
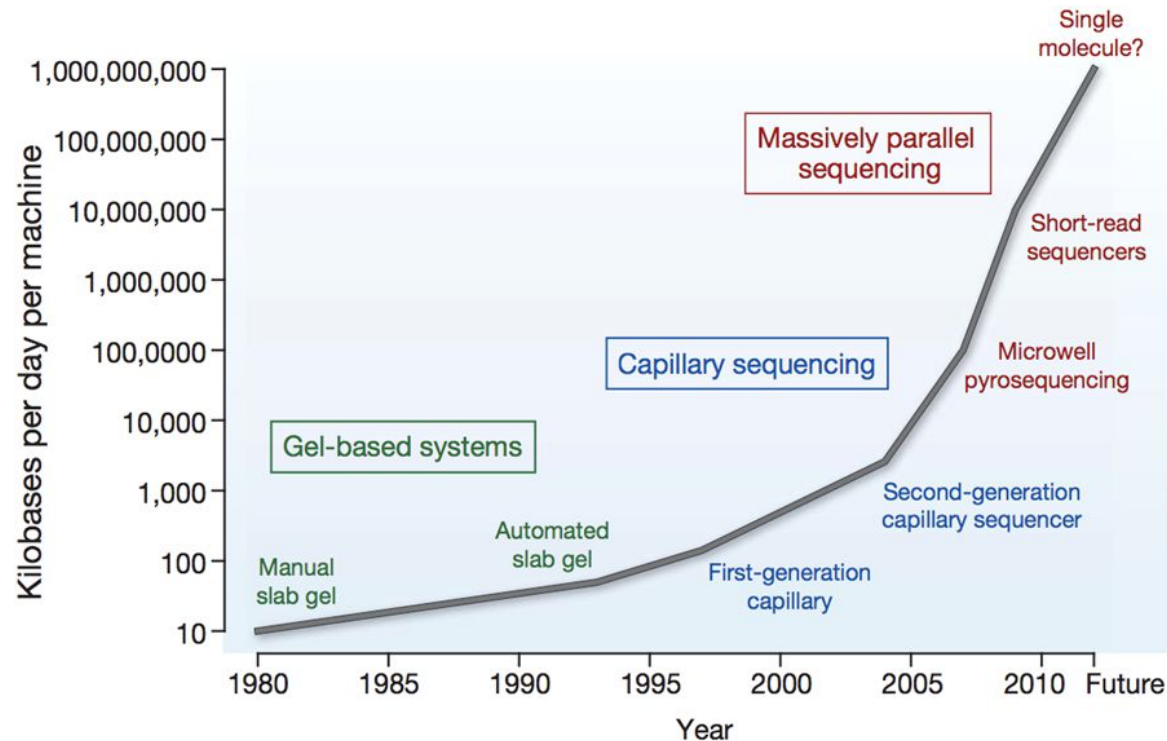
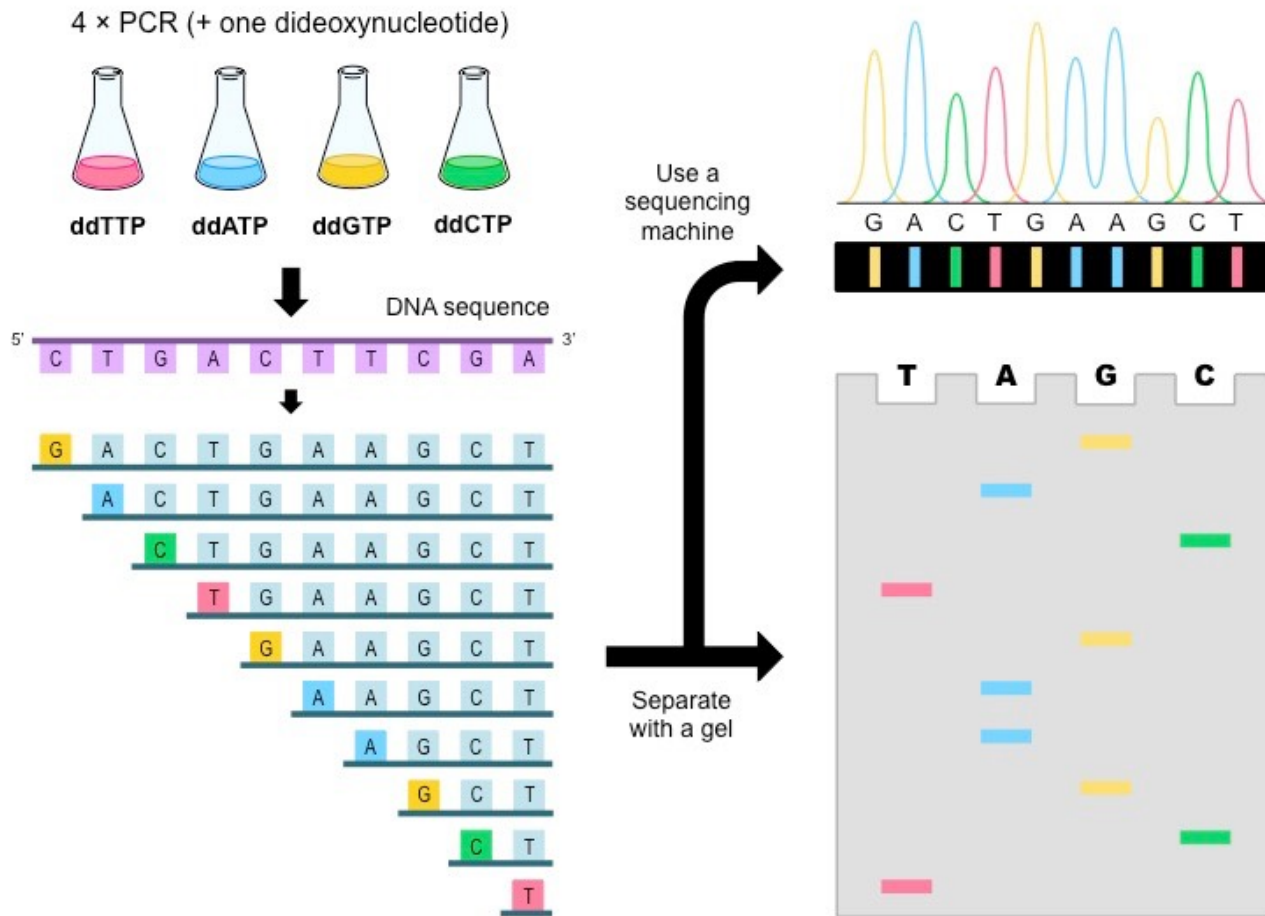# Information including in sequence



- DNA sequence includes not only template for transcription, but also act as regulation factors
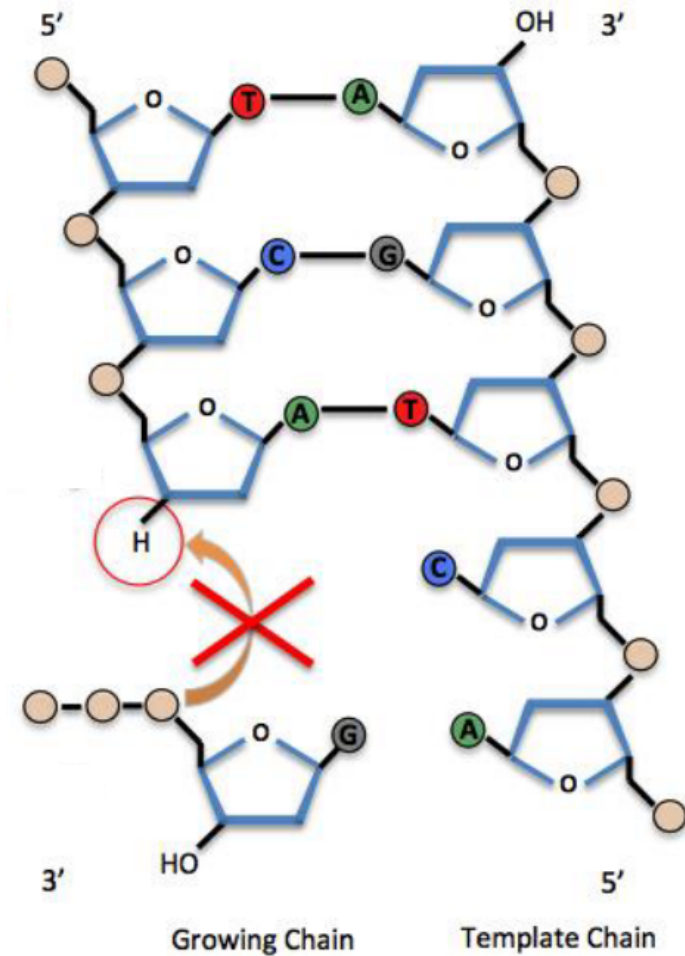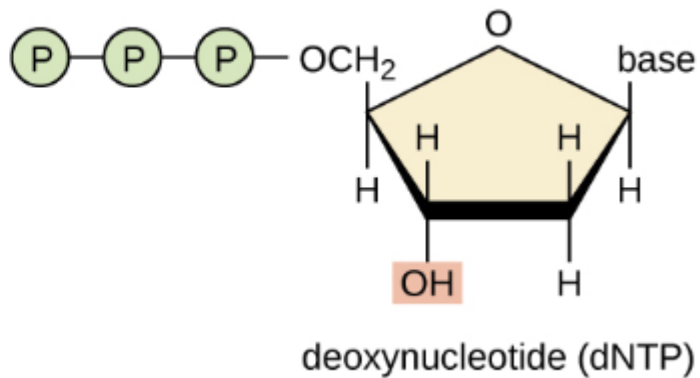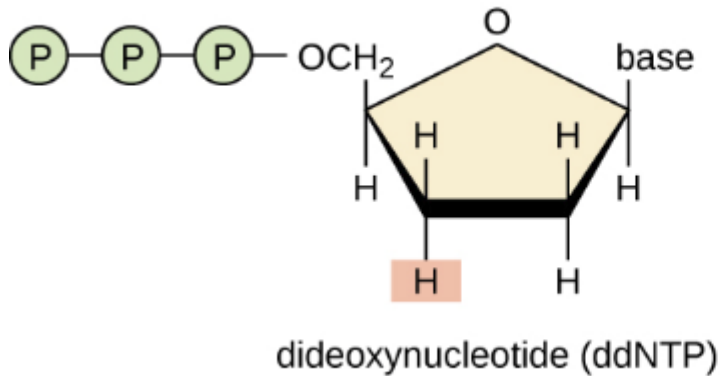
# History of Sequencing technology

# Sanger Sequencing

BMS Bio-Medical Science Co., Ltd.

# Sanger Sequencing

# Sanger Sequencing

| Technology | Analysis time | Average read length | Throughput (Mb/h) |
|---|---|---|---|
| Slab gel | 6–8 hours | 700 bp | 0.0672 |
| Capillary array electrophoresis | 1–3 hours | 700 bp | 0.166 |

**https://en.wikipedia.org/wiki/Sanger_sequencing**

BMS　Bio-Medical Science Co., Ltd.

# Next generation Sequencing

# Next generation Sequencing

# Definition of "Massive parallel"

- **Massive**
  **[형용사]** (육중하면서) 거대한
- **Parallel**
  **[형용사] 병행[병렬]의**

- **Massive parallel sequencing**
  **The DNA is sequenced via spatially separated, clonally amplified DNA templates or single DNA molecules in a flow cell.**

**Massive parallel sequencing = Next Generation Sequencing**

# Next Generation Sequencing Platform

- Roche/454 (GS FLX+/GS Junior)
- Illumina Genome Analyzer (HiSeq/MiSeq/NextSeq)
- Life Technologies (3500 Genetic Analyzer)
- Ion Torrent Proton/PGM)
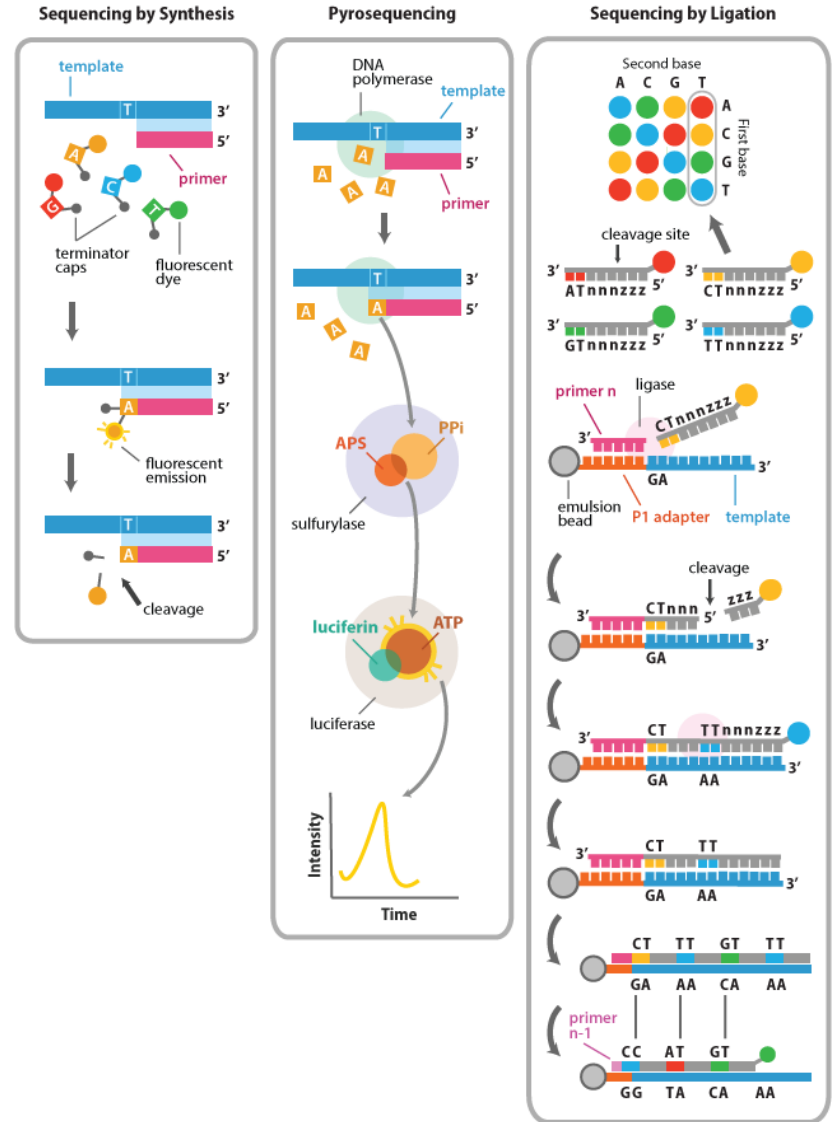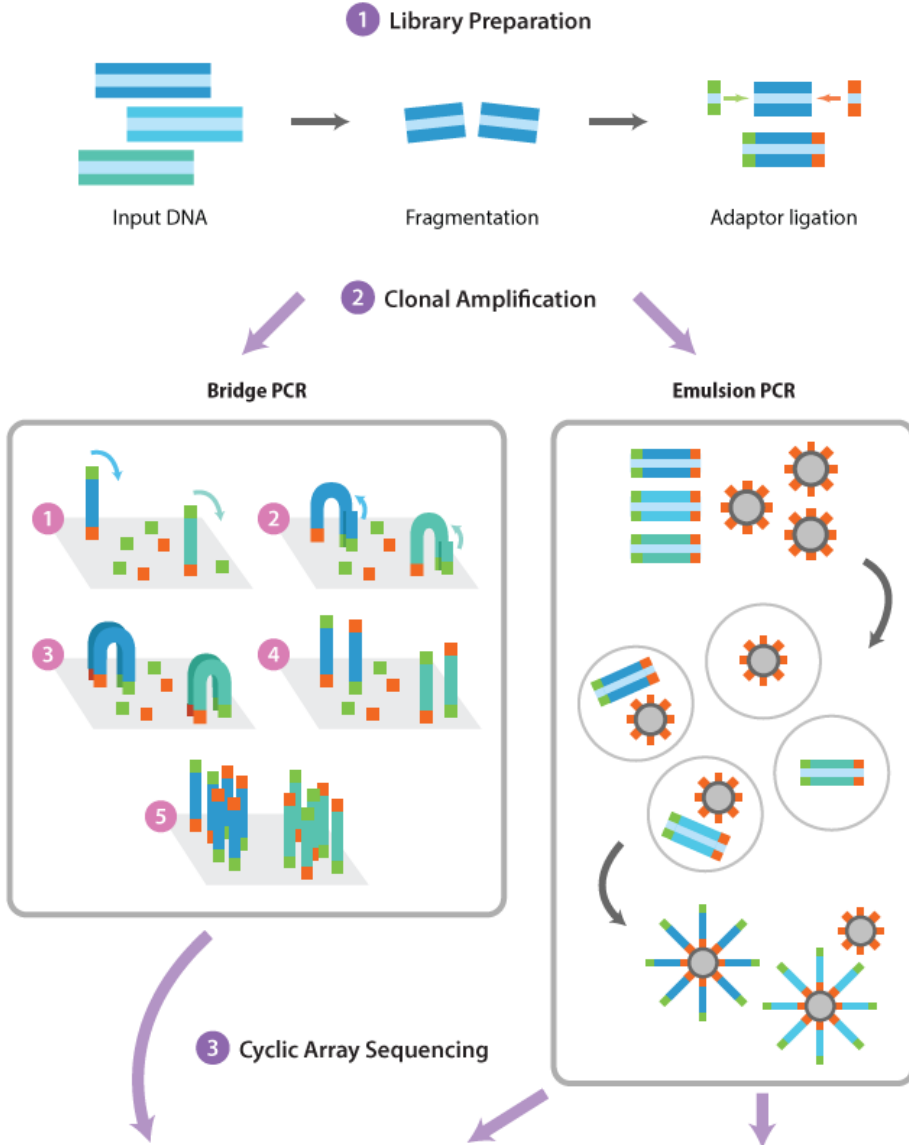- Applied Biosystems (SOLiD, 3730xl DNA Analyzer )

*http://www.vib.be/en/about-vib/annual-report/2012/research/activities/Pages/Service%20Facilities.aspx*
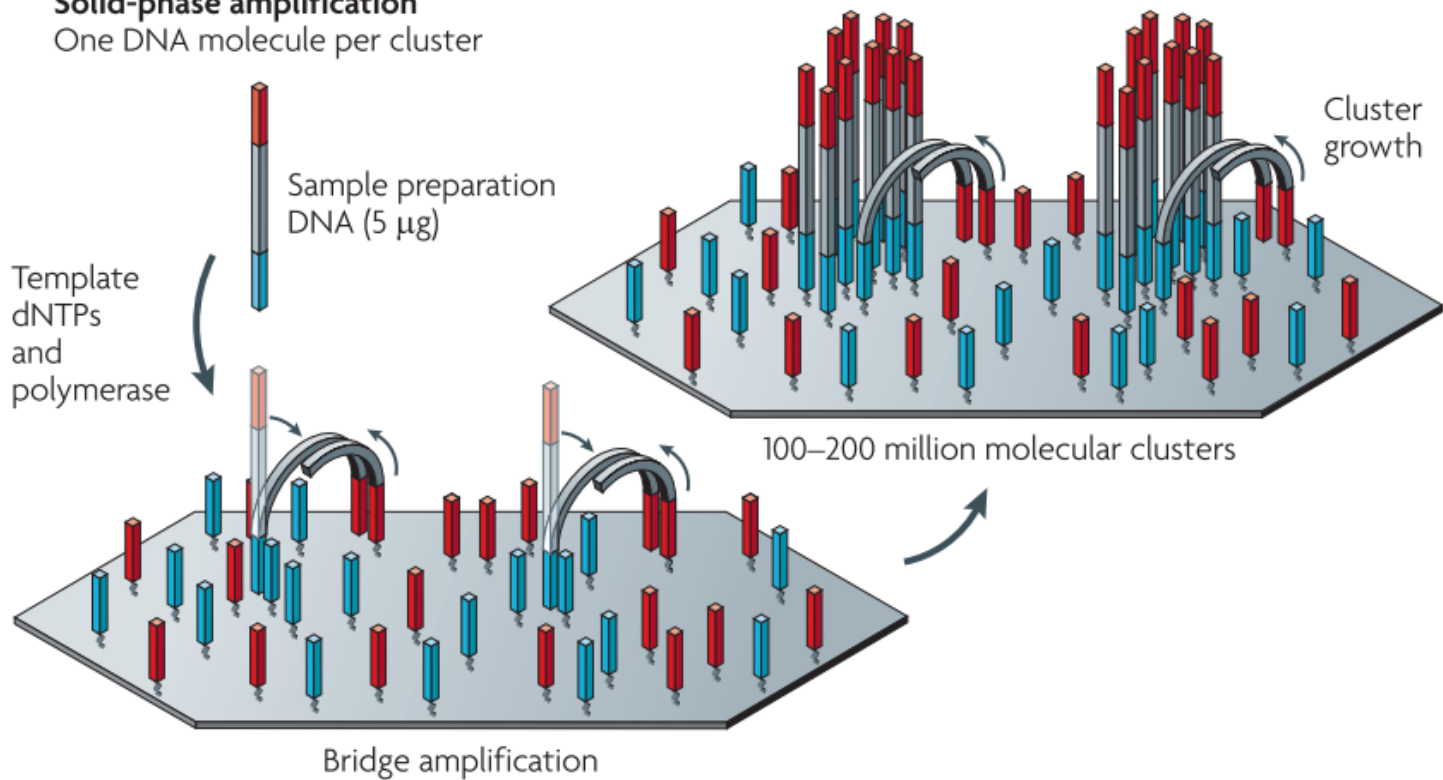
# Next Generation Sequencing Machine

BMS  Bio-Medical Science Co., Ltd.

# Steps of NGS

# Clonal amplification - Bridge PCR



b Illumina/Solexa
Solid-phase amplification
One DNA molecule per cluster

Sample preparation
DNA (5 µg)

Template
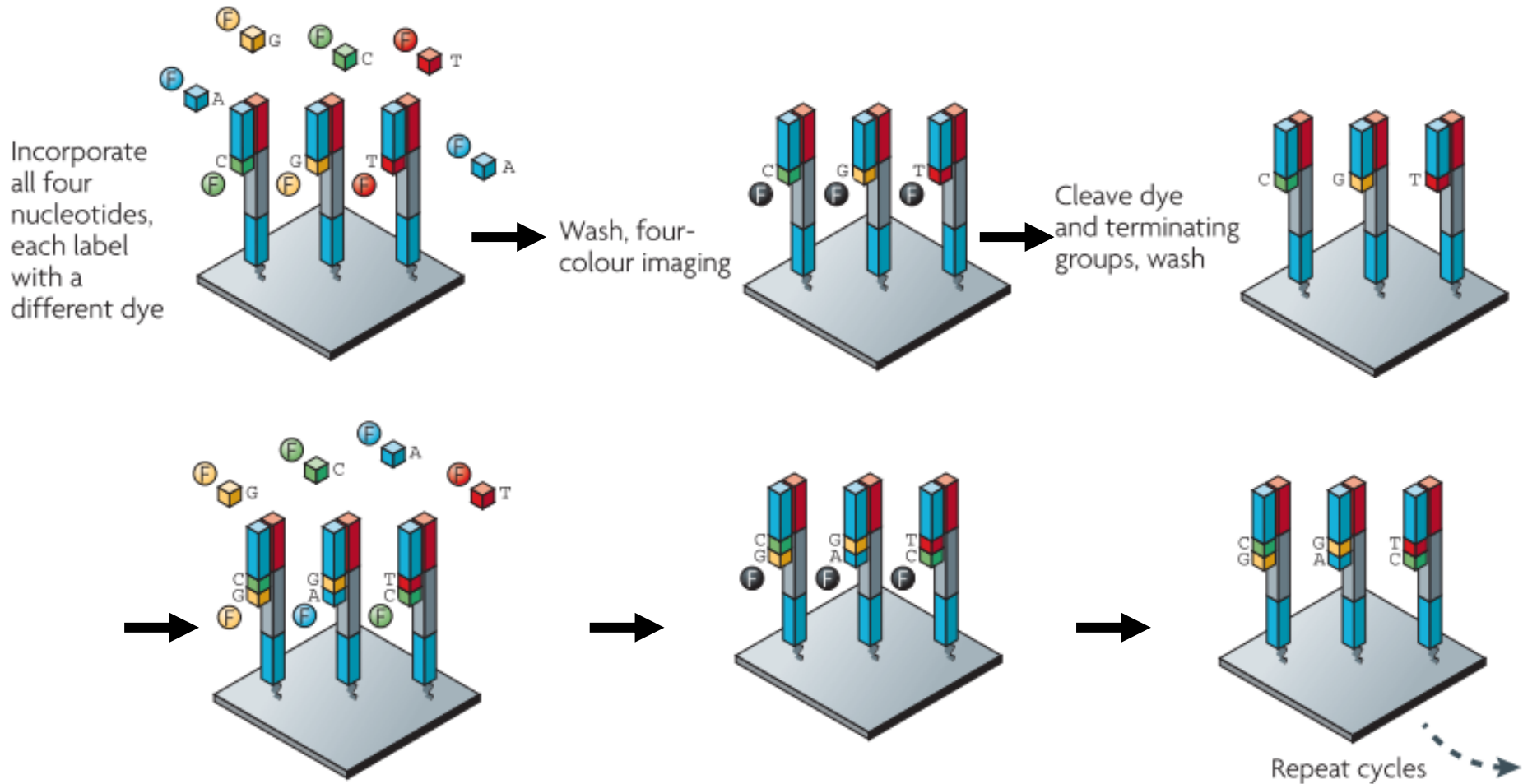dNTPs
and
polymerase

Cluster
growth

100–200 million molecular clusters

Bridge amplification

BMS  Bio-Medical Science Co., Ltd.

# Clonal amplification – Emulsion PCR



**a** Roche/454, Life/APG, Polonator
**Emulsion PCR**
One DNA molecule per bead. Clonal amplification to thousands of copies occurs in microreactors in an emulsion
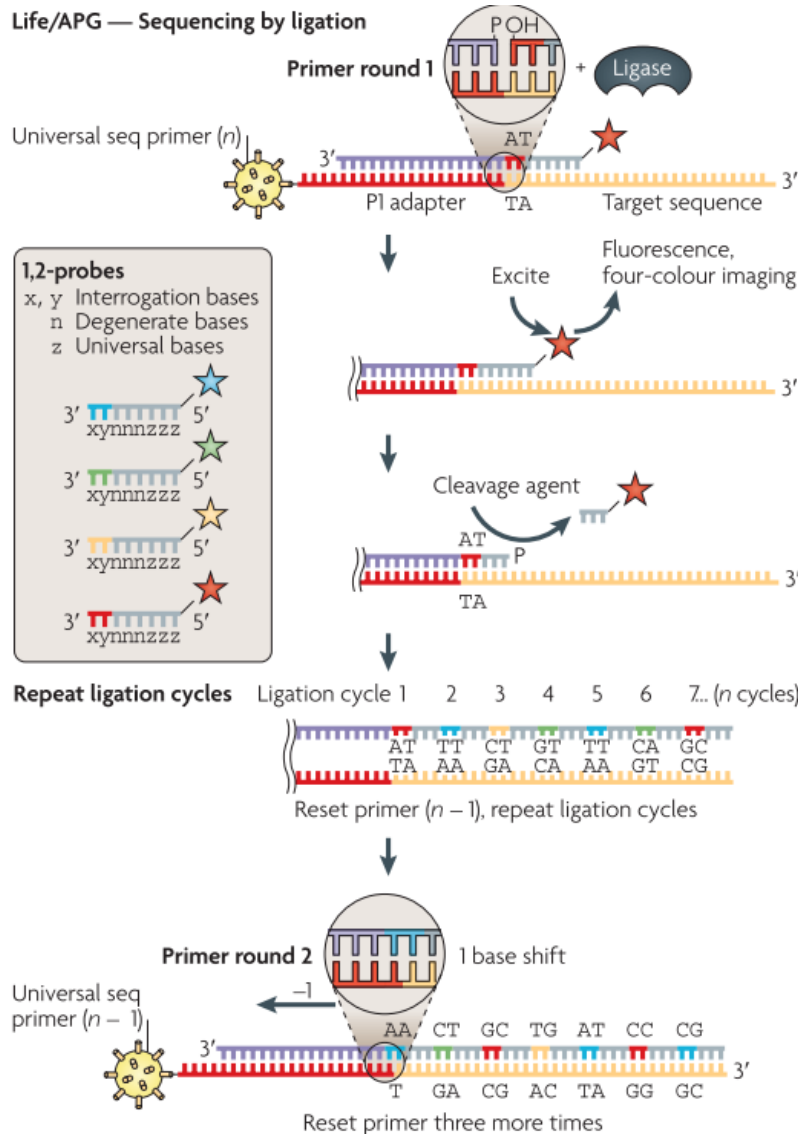
Primer, template, dNTPs and polymerase

PCR amplification

Break emulsion

Template dissociation

100–200 million beads

Chemically cross-linked to a glass slide

# Sequencing by sequencing

# Pyrosequencing



Flowgram

TCAGGTTTTTTAACAATCAACTTTTTGGATTAAAATGTAGATAACTG
CATAAATTAATAACATCACATTAGTCTGATCAGTGAATTTAT

# Sequencing by ligation



Life/APG — Sequencing by ligation

# Comparison between Sanger and NGS


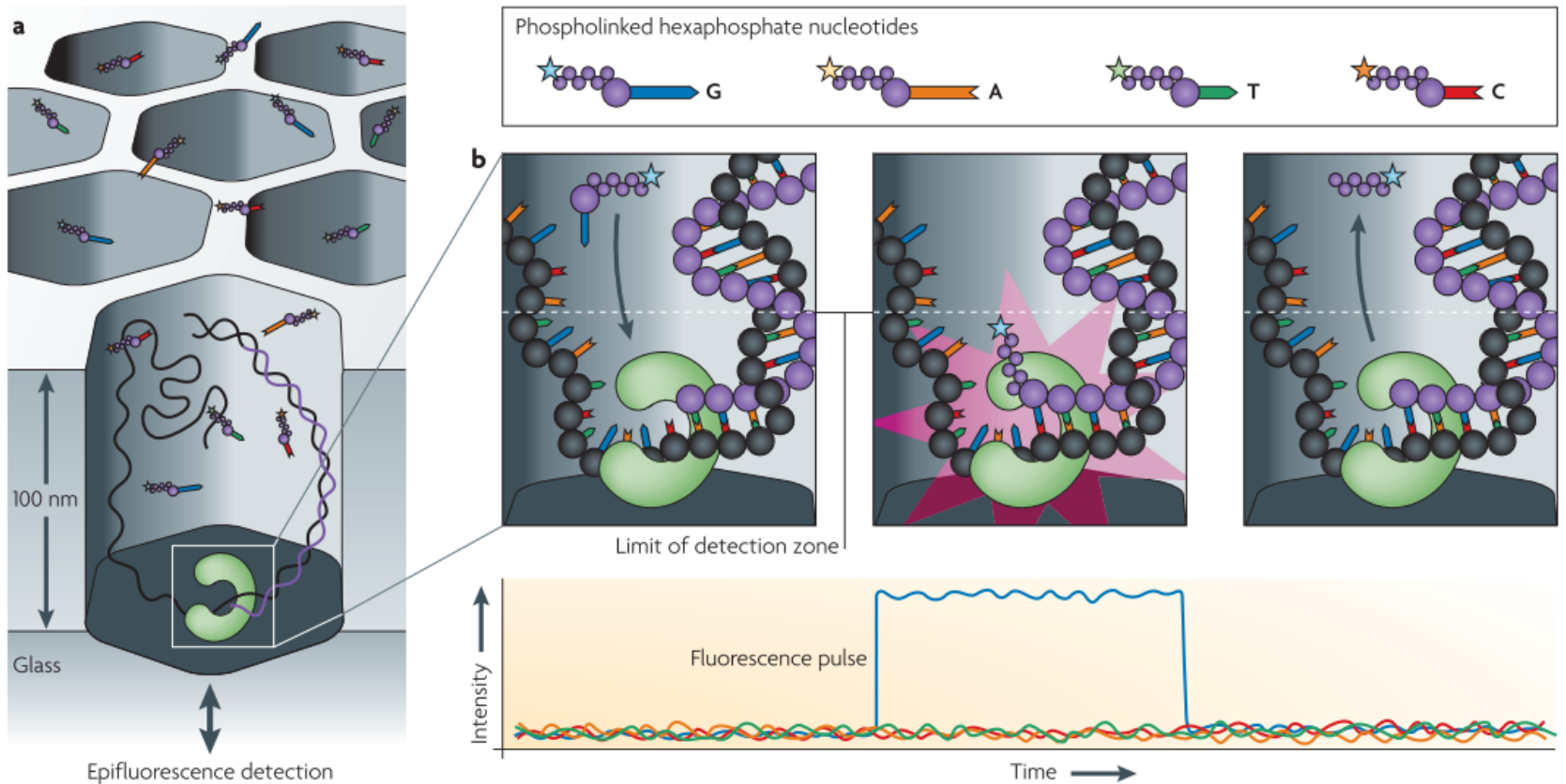
*Nature biotechnology*, 2008

# History of Sequencing technology

# Third Generation Sequencing

- Third generation sequencing works by reading the nucleotide sequences at the single molecule level.

- Existing methods require breaking long strands of DNA into small segments then inferring nucleotide sequences by amplification and synthesis.

- Third Generation Sequencing = long-read sequencing

- PacBio and Oxford Nanopore provide manufactures

*Nature biotechnology*, 2008

# PacBio

# Steps in PacBio



1. generate amplicon

5' forward strand 3'
3' reverse strand 5'

2. ligate adaptors

SMRTbell

3. sequence

template
DNA polymerase

4. data analysis

raw long read
processed long read
single-molecule fragments
circular consensus sequence (ccs)

1° analysis

# Nanopore

# Steps in Nanopre



(i)  (ii)  (iii)  (iv)  (v)  (vi)  (vii)  (viii)

**BMS** Bio-Medical Science Co., Ltd.

# Strategy to improve sequencing accuracy

# Accuracy of Nanopore

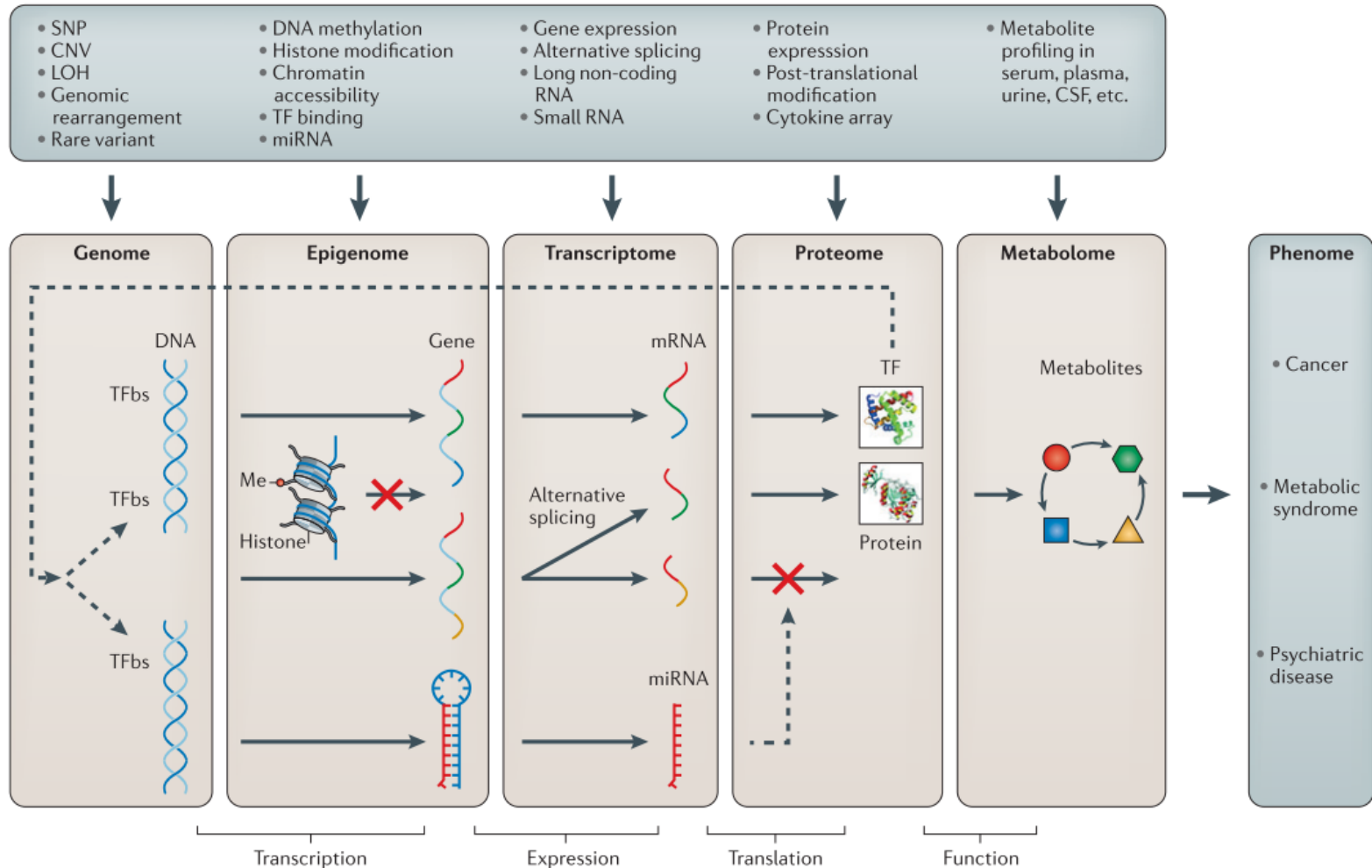BMS  Bio-Medical Science Co., Ltd.

# Comparison of platform between 2<sup>nd</sup> and 3<sup>rd</sup>

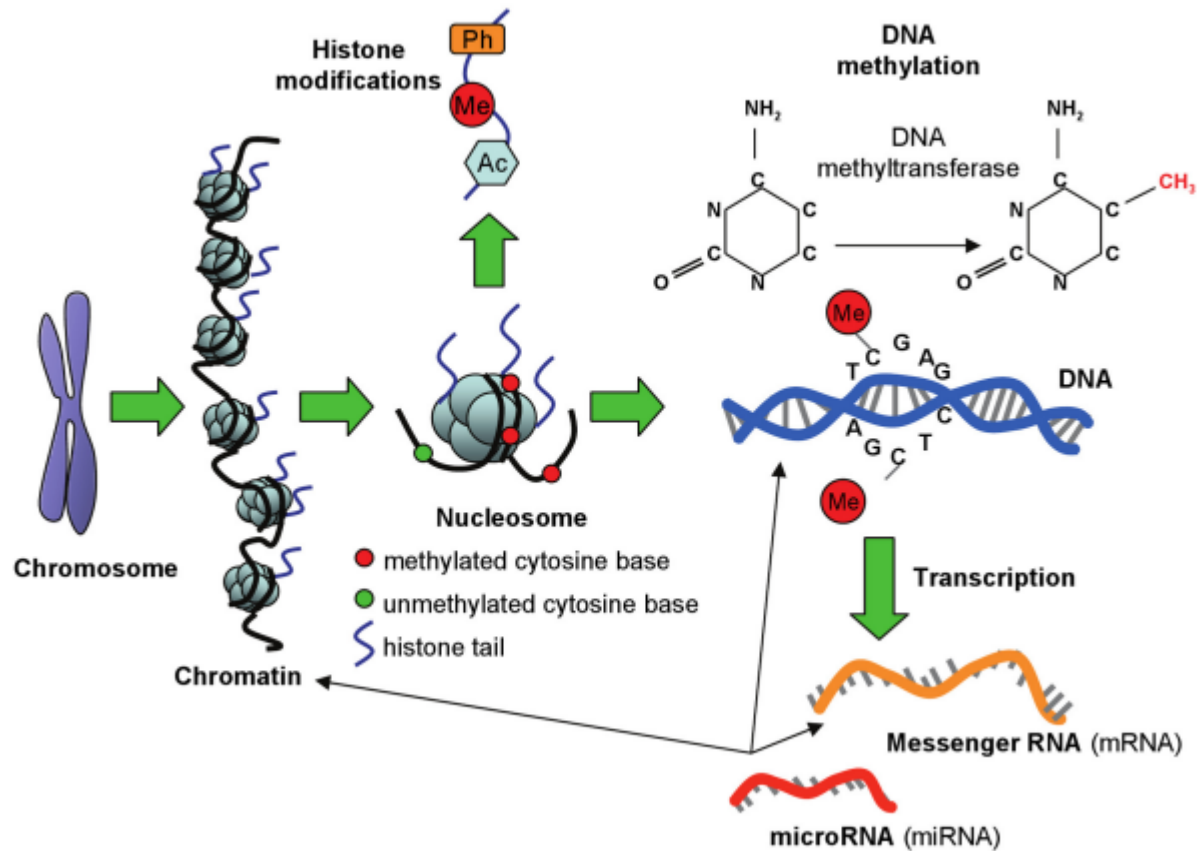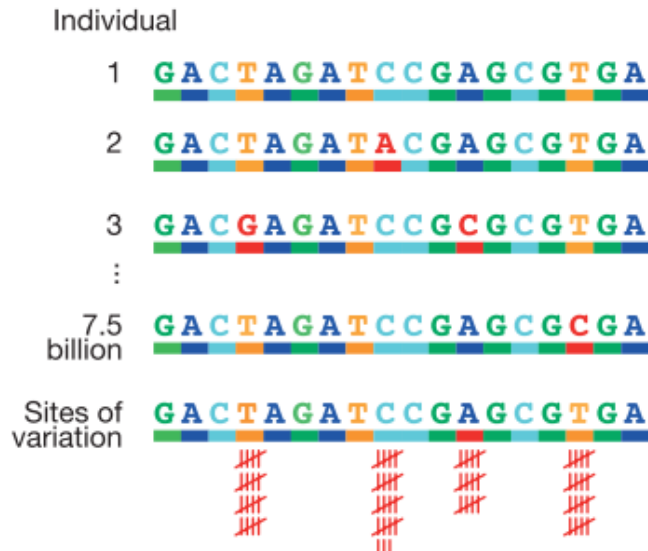| Sequencing platform | | Total output (bases per run) | Total reads (million per run) | Read length (bases) | Run time (days) | Purpose/definition |
|---|---|---|---|---|---|---|
| Illumina | HiSeq X | 1.6−1.8 Tb | 6000 M | 2 × 150 bp | <3 | Allows sequencing of larger genomes (e.g., mammalian genomes) at population level |
| | MiSeq | 300 Mb−15 Gb | 50 M | 2 × 300 bp | 0.2−2.7 | Designed for particularly small genomes (e.g., bacterial genomes) and amplicon sequencing |
| Life technologies | Solid 5500 Systems | 80 Gb−320 Gb | 1200 M−2400 M | 50−2 × 50 bp | 7 | Offers application-per-lane sequencing that allows transcriptome, exome and genome sequencing concurrently in a single run. Additionally, pay-per-lane sequencing feature makes Solid 5500 Systems cost-effective because reagents are not required for unused lanes. |
| | Ion Torrent 520 Chip | 600 Mb−2 Gb | 3−5 M | 200−400 bp | 0.1 | Ion S5 System allows generation of diverse sequencing data ranging from targeted re-sequencing to genome sequencing with as little as 10 ng sample. |
| | Ion Torrent 540 Chip | 10−15 Gb | 60−80 M | 200−400 bp | 0.1 | |
| PacBio | Sequel System | 500 Mb−16 Gb | 55−880 M | up to 60 kb | <0.1−0.3 | Useful in the studies of *de novo* assembly of large genomes. Sequel System can be utilized for generating variation, expression and/or regulation related sequencing data. |
| | PacBio RS II | 500 Mb−16 Gb | 55−880 M | up to 60 kb | <0.1−0.3 | Much more suitable for sequencing small genomes although animal and plant genomic studies is also possible. |
| Nanopore | PromethION | up to 12 Tb[a] | 1250 M[a] | 230−300 kb[a] | 2 | Ideal for large sample numbers. PromethION can sequence up to 48 samples in a single run |
| | MinION | up to 42 Gb[a] | up to 4.4 M[a] | 230−300 kb[a] | 2 | Portable sequencing instrument. MinION can be run with a desktop or laptop computer and data can be performed in real time. |

# Application of NGS
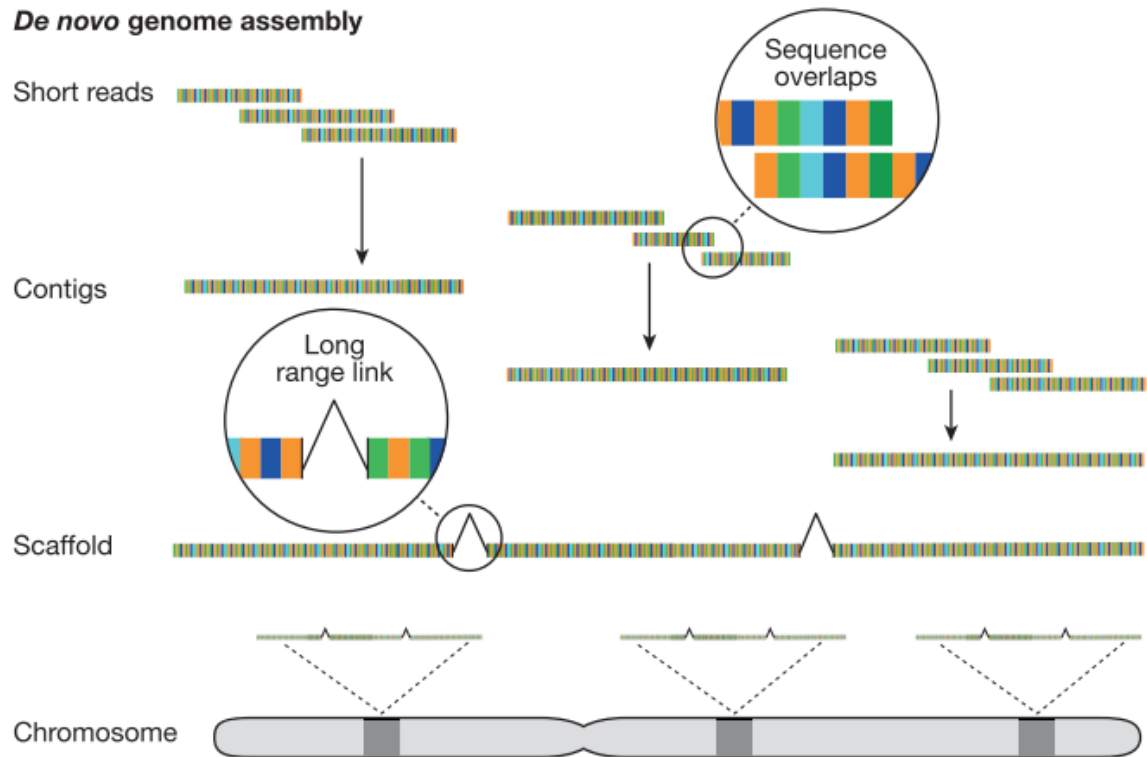
# Biological layer

# Central dogma with epigenome

BMS    Bio-Medical Science Co., Ltd.

# Sequencing using whole genome

BMS  Bio-Medical Science Co., Ltd.

# Type of variants



**Types of Variants**

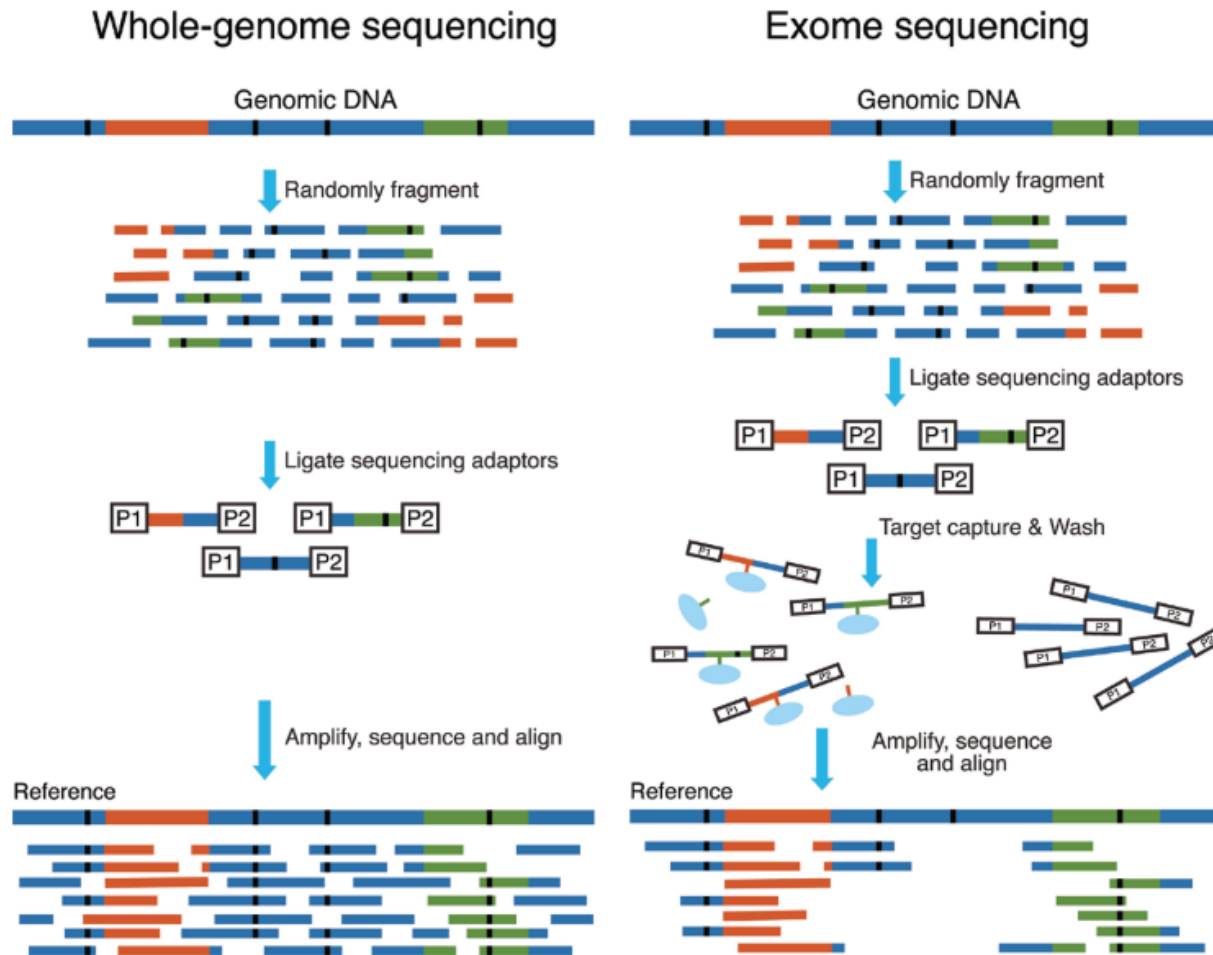*https://www.pacb.com/applications/whole-genome-sequencing/variant-detection/*

# The Angelina effect

# BRCA1 and BRCA2

*BRCA1* and *BRCA2* are human genes that produce tumor suppressor proteins. These proteins help repair damaged DNA and, therefore, play a role in ensuring the stability of each cell's genetic material. When either of these genes is mutated, or altered, such that its protein product is not made or does not function correctly, DNA damage may not be repaired properly. As a result, cells are more likely to develop additional genetic alterations that can lead to cancer.

Specific inherited mutations in *BRCA1* and *BRCA2* most notably increase the risk of female breast and ovarian cancers, but they have also been associated with increased risks of several additional types of cancer. People who have inherited mutations in *BRCA1* and *BRCA2* tend to develop breast and ovarian cancers at younger ages than people who do not have these mutations.

**Breast cancer:** About 12% of women in the general population will develop breast cancer sometime during their lives (1). By contrast, a recent large study estimated that about 72% of women who inherit a harmful *BRCA1* mutation and about 69% of women who inherit a harmful *BRCA2* mutation will develop breast cancer by the age of 80 (2).

**Ovarian cancer:** About 1.3% of women in the general population will develop ovarian cancer sometime during their lives (1). By contrast, it is estimated that about 44% of women who inherit a harmful *BRCA1* mutation and about 17% of women who inherit a harmful *BRCA2* mutation will develop ovarian cancer by the age of 80 (2).

**https://www.cancer.gov/about-cancer/causes-prevention/genetics/brca-fact-sheet**

# Sequencing for resequencing

BMS  Bio-Medical Science Co., Ltd.

# Whole Exome Sequencing

BMS  Bio-Medical Science Co., Ltd.

# Steps of De novo genome assembly

BMS    Bio-Medical Science Co., Ltd.

# Genome assembly algorithms

## A Read Layout

R₁: GACCTACA
R₂:   ACCTACAA
R₃:     CCTACAAG
R₄:       CTACAAGT
A:          TACAAGTT
B:            ACAAGTTA
C:              CAAGTTAG
X:          TACAAGTC
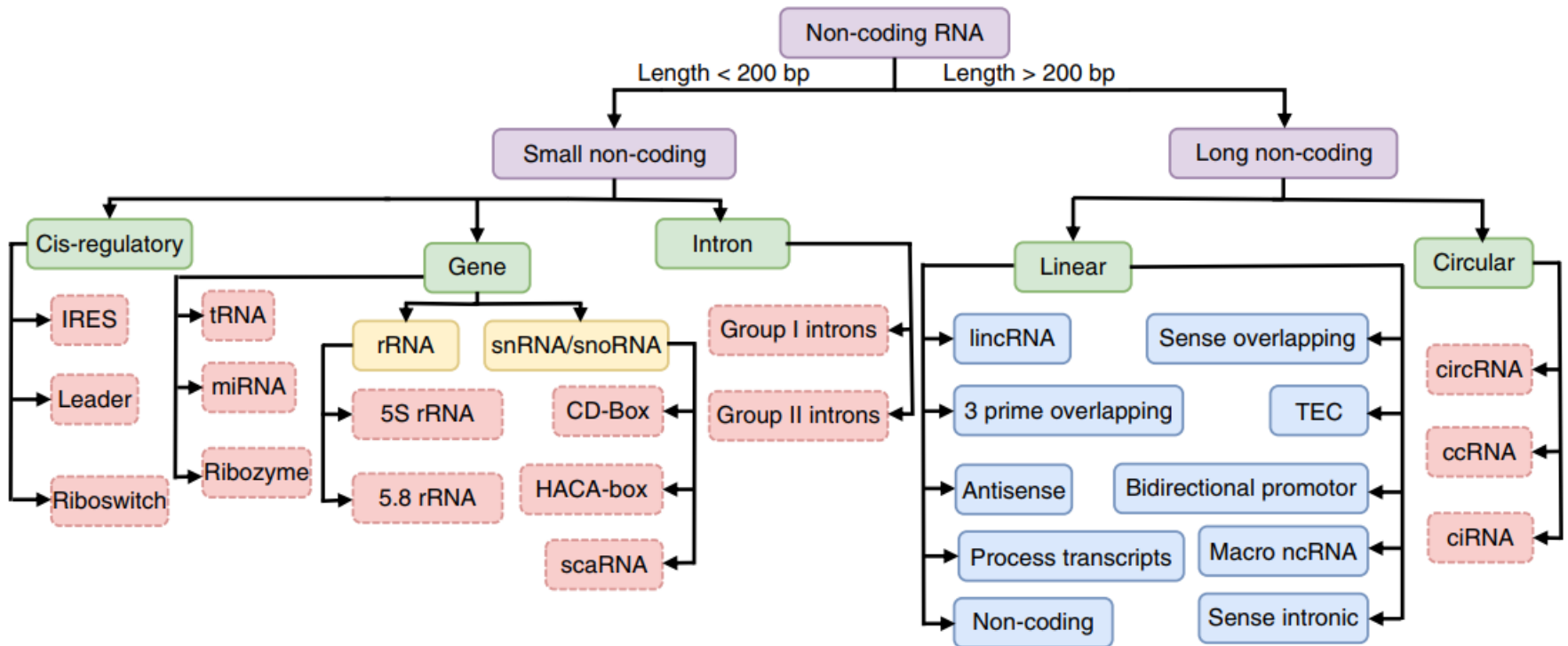Y:            ACAAGTCC
Z:              CAAGTCCG

## B Overlap Graph

## C de Bruijn Graph

# Messenger RNA



Prokaryotic mRNA:

Eukaryotic mRNA:

BMS    Bio-Medical Science Co., Ltd.

# Non-coding RNA

# Type of RNA-seq

# Type of RNA-seq

**BMS** Bio-Medical Science Co., Ltd.

# Type of RNA-seq

# RNA-seq data analysis

**BMS** Bio-Medical Science Co., Ltd.

# RNA-seq data analysis



*Nature reviews genetics, 2009*

# Nobel prize – RNA Interference

The Nobel Prize in Physiology or Medicine 2006

Photo: L. Cicero
Andrew Z. Fire
Prize share: 1/2

Photo: J. Mottern
Craig C. Mello
Prize share: 1/2

*https://www.nobelprize.org/prizes/medicine/2006/summary/*

BMS  Bio-Medical Science Co., Ltd.

# Micro RNA biogenesis



*Nature reviews Genetics, 2004*

# Chromatin and Chromatin immunoprecipitation

# Sequencing Application in Epigenome

# Epigenome function

# ChIP-Seq and DNA accessibility sequencing

# Histone modification



*Nature reviews Cancer, 2001*

# Regulation by histone profile

# Active regulation in gene

# 'Dashboard' of histone modifications



*Nature reviews. Genetics, 2011*

# Histone profile depending on Cell type

BMS  Bio-Medical Science Co., Ltd.

# Classes of DNA-bound proteins

# Principle detection of protein binding region

# Principle detection of protein binding region

Design and analysis of ChIP-seq experiments for
DNA-binding proteins

Peter V Kharchenko[1-3], Michael Y Tolstorukov[1,2] & Peter J Park[1-3]

An integrated software system for analyzing ChIP-chip
and ChIP-seq data

Hongkai Ji[1], Hui Jiang[2], Wenxiu Ma[3], David S Johnson[4,8], Richard M Myers[5] & Wing H Wong[6,7]

## FindPeaks 3.1: a tool for identifying areas of enrichment from massively parallel short-read sequencing technology

Anthony P. Fejes[1,*], Gordon Robertson[1], Mikhail Bilenky[1], Richard Varhol[1],
Matthew Bainbridge[2] and Steven J. M. Jones[1,*]

HPeak: an HMM-based algorithm for defining
read-enriched regions in ChIP-Seq data

Zhaohui S Qin[*1,2,3], Jianjun Yu[3,4], Jincheng Shen[1], Christopher A Maher[2,3,4], Ming Hu[1], Shanker Kalyana-Sundaram[3,4]
Jindan Yu[5] and Arul M Chinnaiyan[2,3,4,6,7,8]

## Model-based Analysis of ChIP-Seq (MACS)

Yong Zhang[¤*], Tao Liu[¤*], Clifford A Meyer[*], Jérôme Eeckhoute[†]
David S Johnson[‡], Bradley E Bernstein[§¶], Chad Nusbaum[¶],
Richard M Myers[¥], Myles Brown[†], Wei Li[#] and X Shirley Liu[*]

PeakSeq enables systematic scoring of ChIP-seq
experiments relative to controls

Joel Rozowsky[1], Ghia Euskirchen[2], Raymond K Auerbach[3], Zhengdong D Zhang[1], Theodore Gibson[1],
Robert Bjornson[4], Nicholas Carriero[4], Michael Snyder[1,2] & Mark B Gerstein[1,3,4]

## BayesPeak: Bayesian analysis of ChIP-seq data

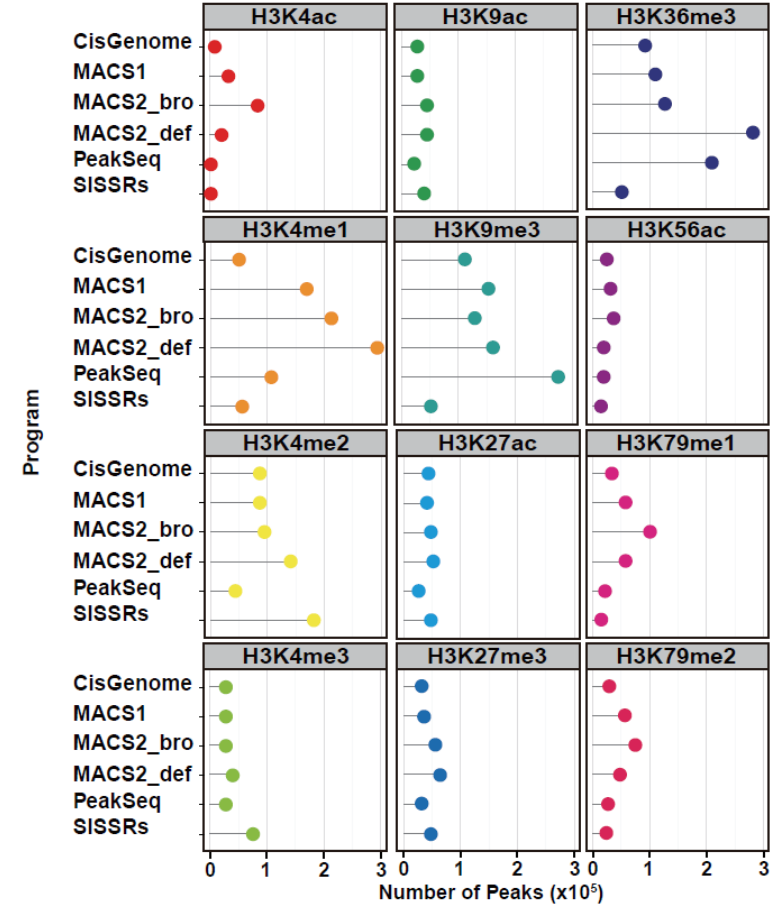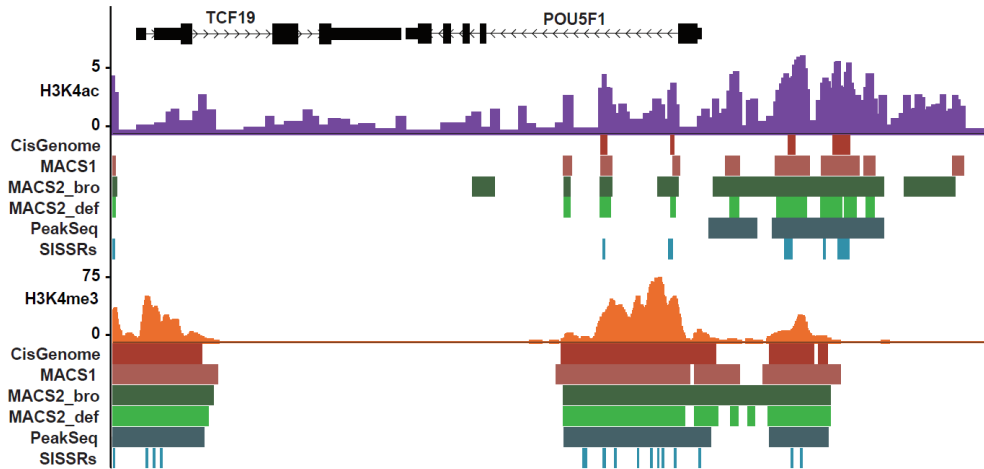Christiana Spyrou[*1,3], Rory Stark[3], Andy G Lynch[4] and Simon Tavaré[2,4]

Genome-wide identification of in vivo protein–DNA
binding sites from ChIP-Seq data

Raja Jothi, Suresh Cuddapah, Artem Barski, Kairong Cui and Keji Zhao[*]

Sole-Search: an integrated analysis program for
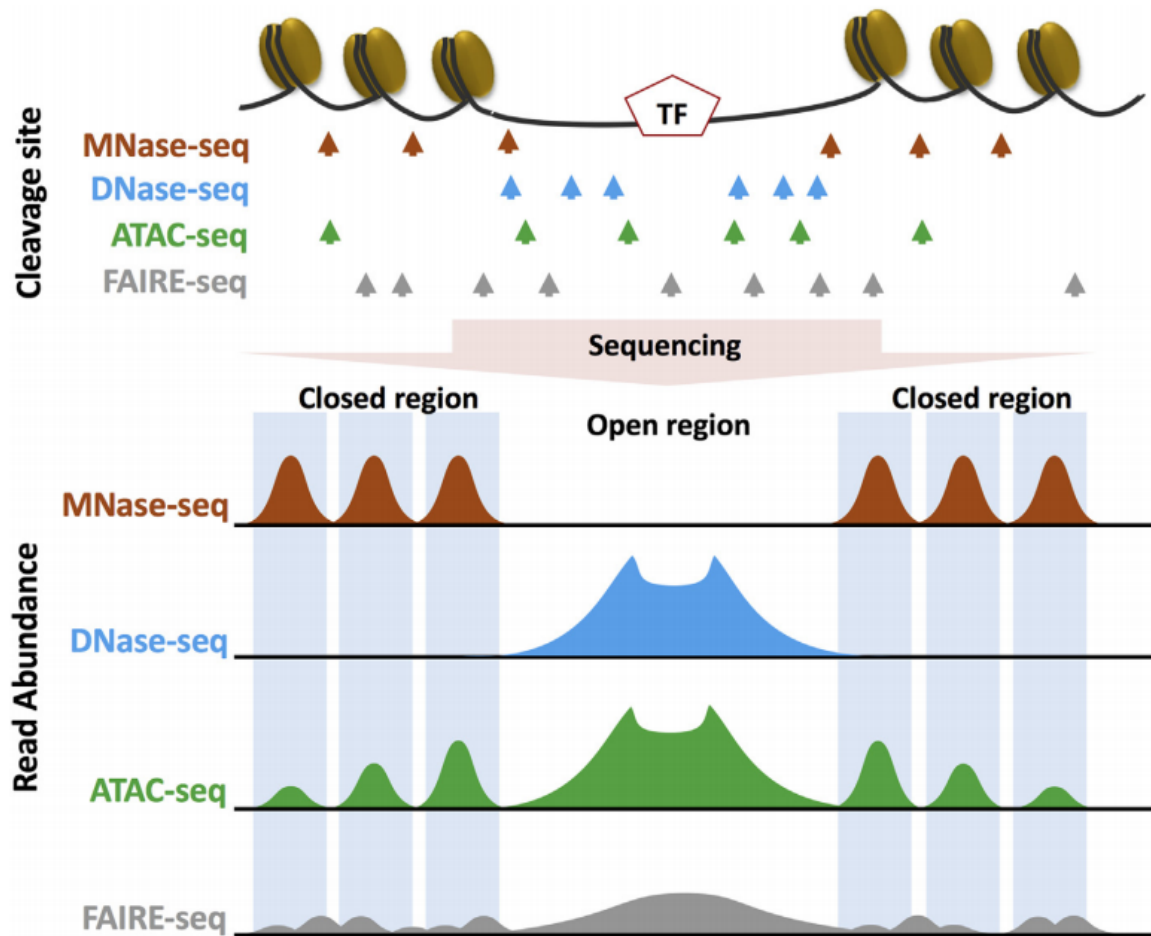peak detection and functional annotation using
ChIP-seq data

Kimberly R. Blahnik[1], Lei Dou[1], Henriette O'Geen[1], Timothy McPhillips[1], Xiaoqin Xu[1],
Alina R. Cao[1], Sushma Iyengar[1], Charles M. Nicolet[1], Bertram Ludäscher[1,2], Ian Korf[1,3]
and Peggy J. Farnham[1,4,*]

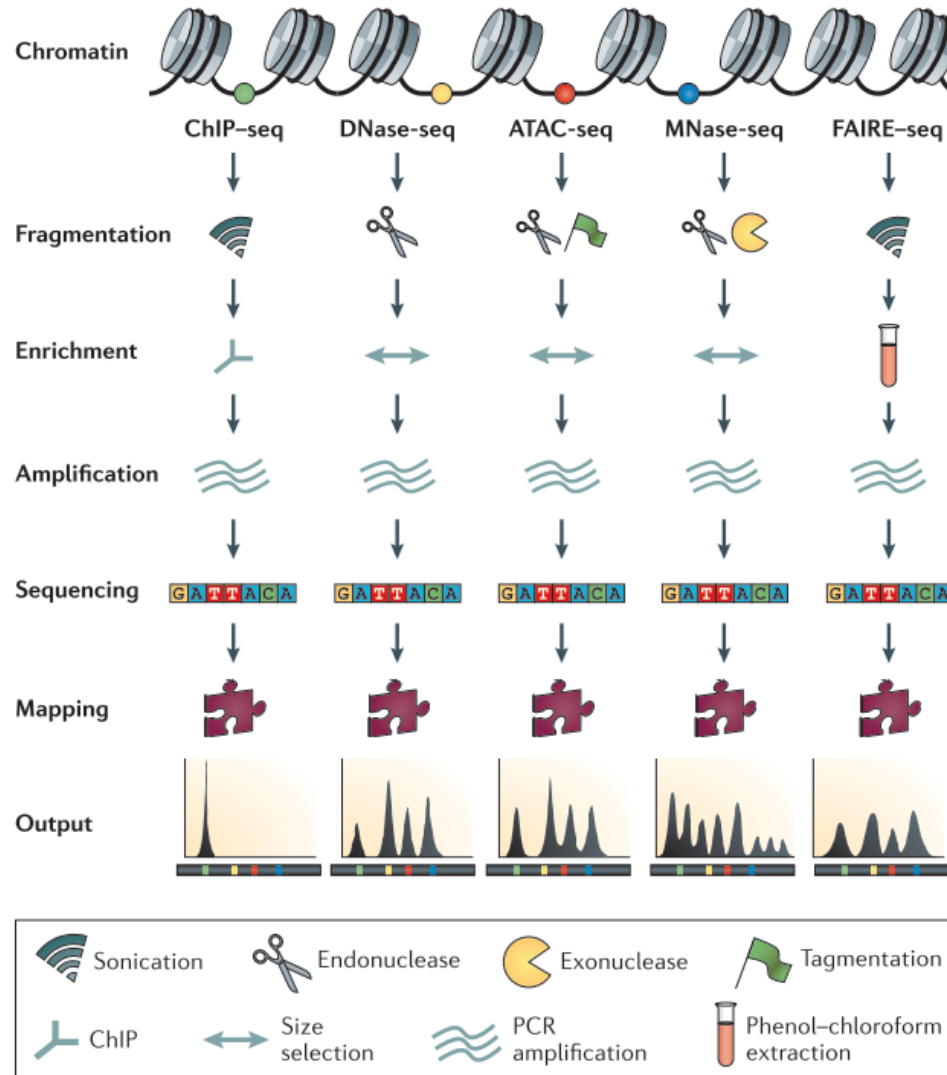# Difference of peak result depending on program

BMS Bio-Medical Science Co., Ltd.

# Active regulation in gene
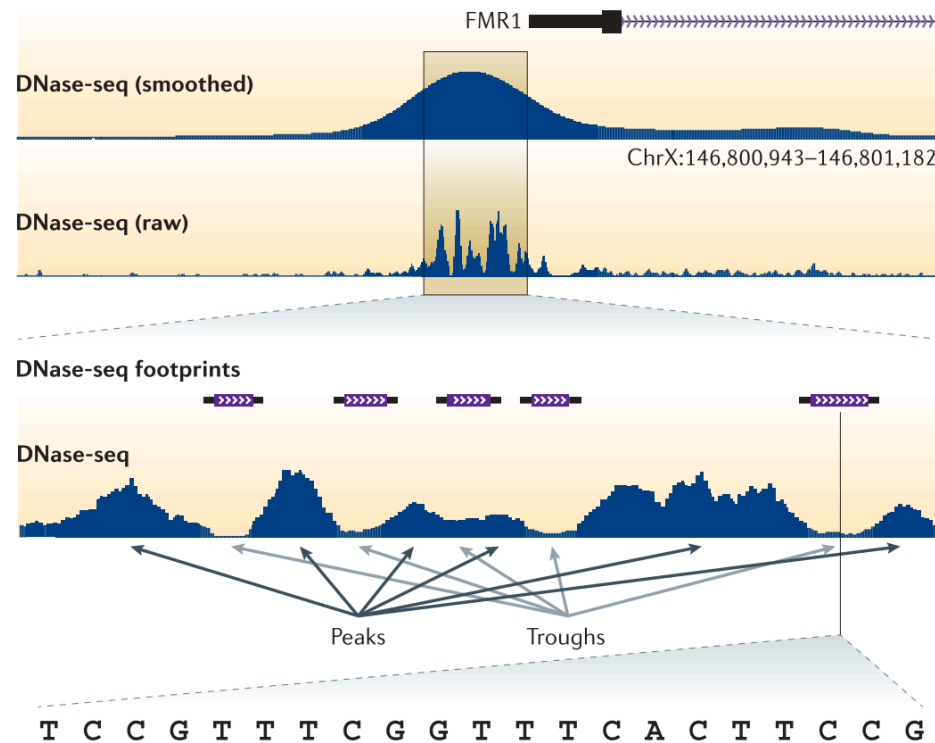
BMS  Bio-Medical Science Co., Ltd.

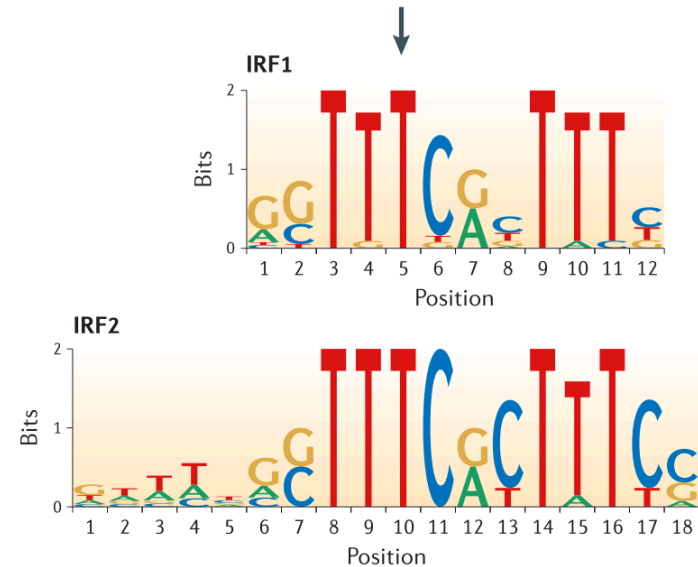# Sequencing to investigate DNA accessibility



*Nature Reviews Genetics, 2014*

# Sequencing to investigate DNA accessibility



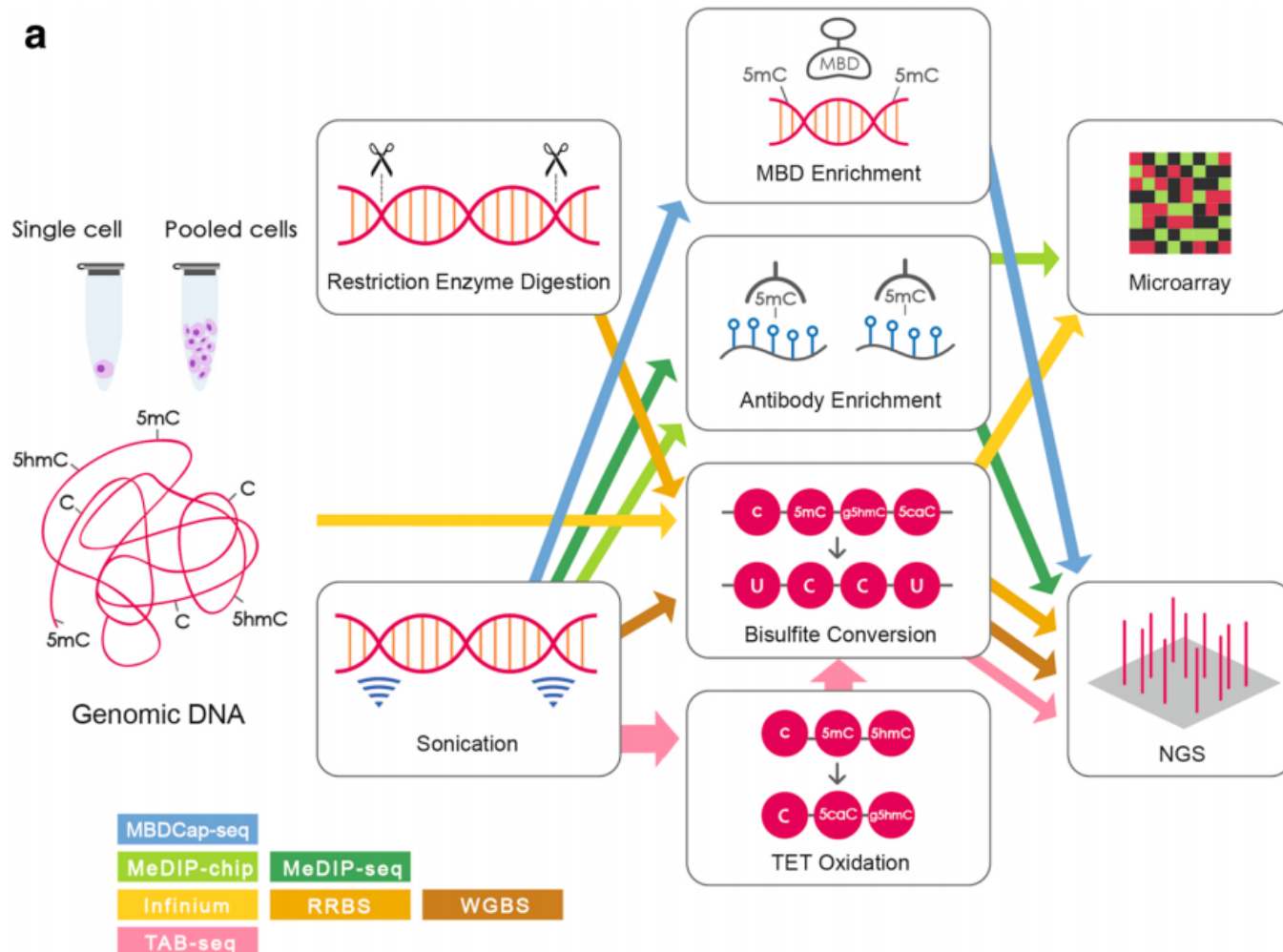Motifs from JASPAR database

| Model name | Score | Relative score | Start | End | Strand | Predicted site sequence |
|---|---|---|---|---|---|---|
| IRF1 | 12.986 | 0.904279917181229 | 3 | 14 | −1 | GAAACCGAAACG |
| IRF2 | 17.216 | 0.907706906384892 | 4 | 21 | −1 | CGGAAGTGAAACCGAAAC |
| SPIB | 4.820 | 0.806987596140569 | 5 | 11 | −1 | ACCGAAA |
| BRCA1 | 4.228 | 0.802287513481405 | 8 | 14 | −1 | GAAACCG |

*Nature Reviews Genetics, 2014*

# Sequencing method to measure methylation



*Cell, 2014*

# DNA methylation mechanism

# DNA methylation mechanism

BMS Bio-Medical Science Co., Ltd.

# Sequencing for methylation

# DNA methylation function

# Chromatin structure

# Chromatin Conformation Capture(3C)

# 3C-based approaches

# Conformation in genome

# Conformation in genome
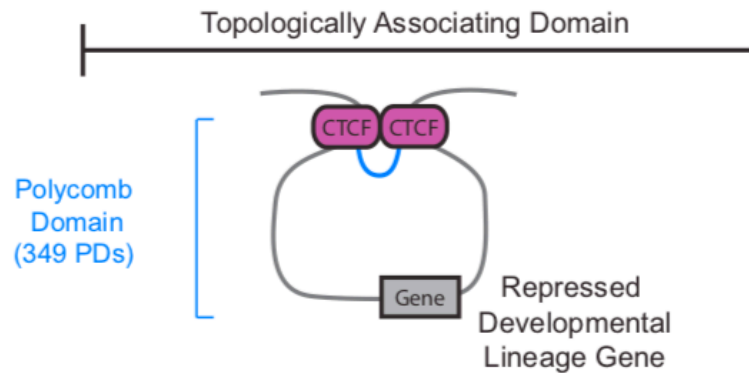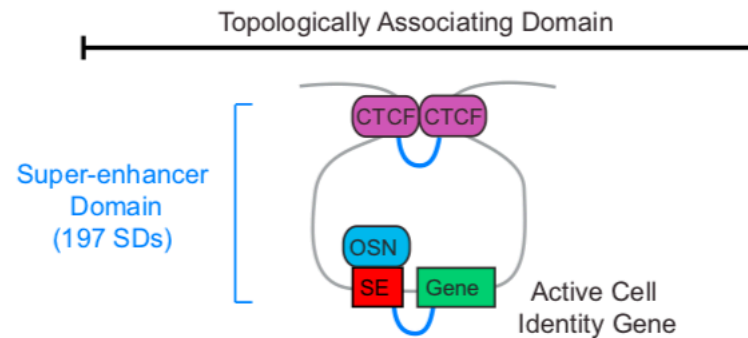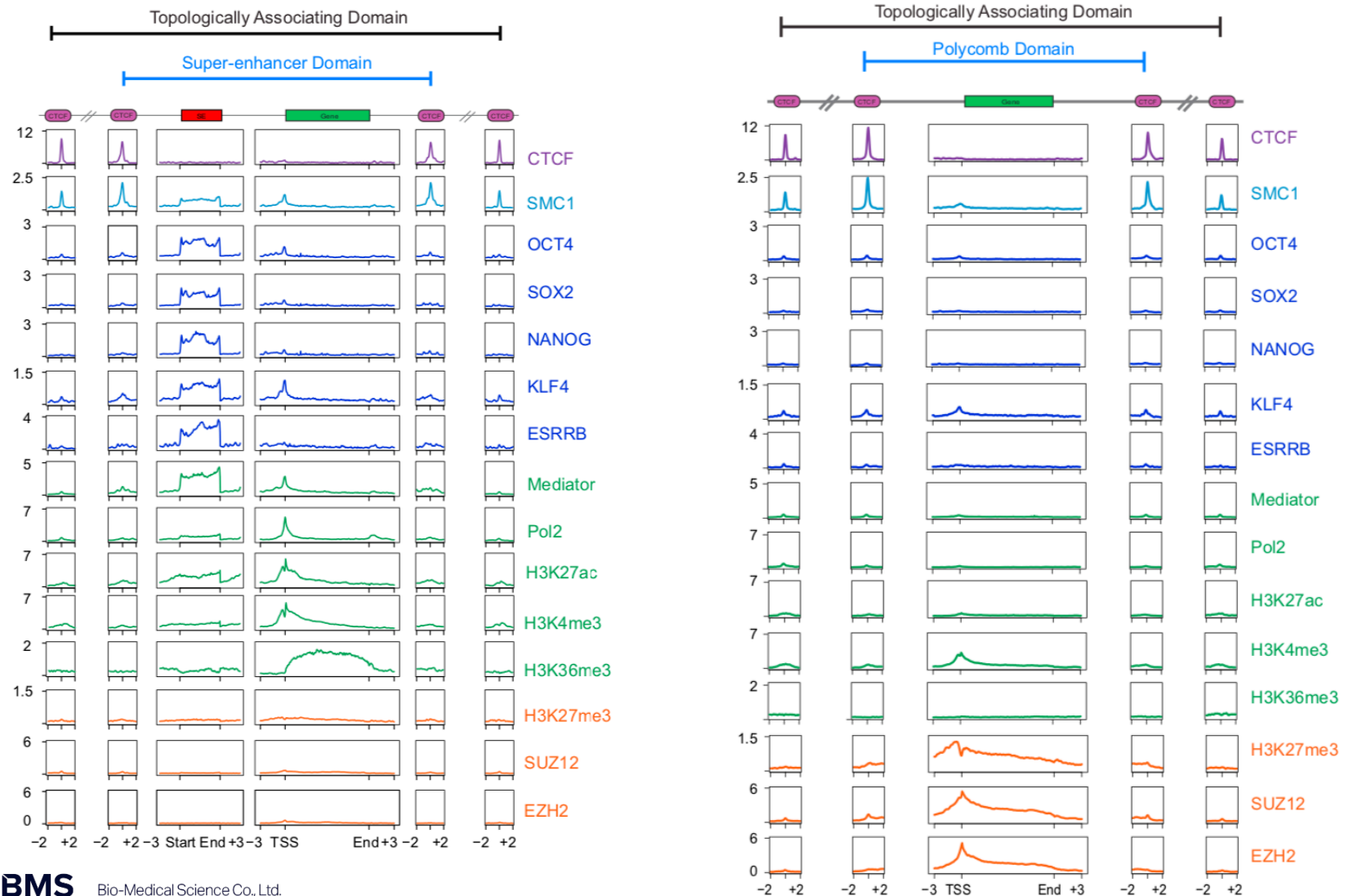
# Diseases associated to chromatin conformation

# Interplay among epigenome feature

# Interplay among epigenome feature

# Interplay among epigenome feature

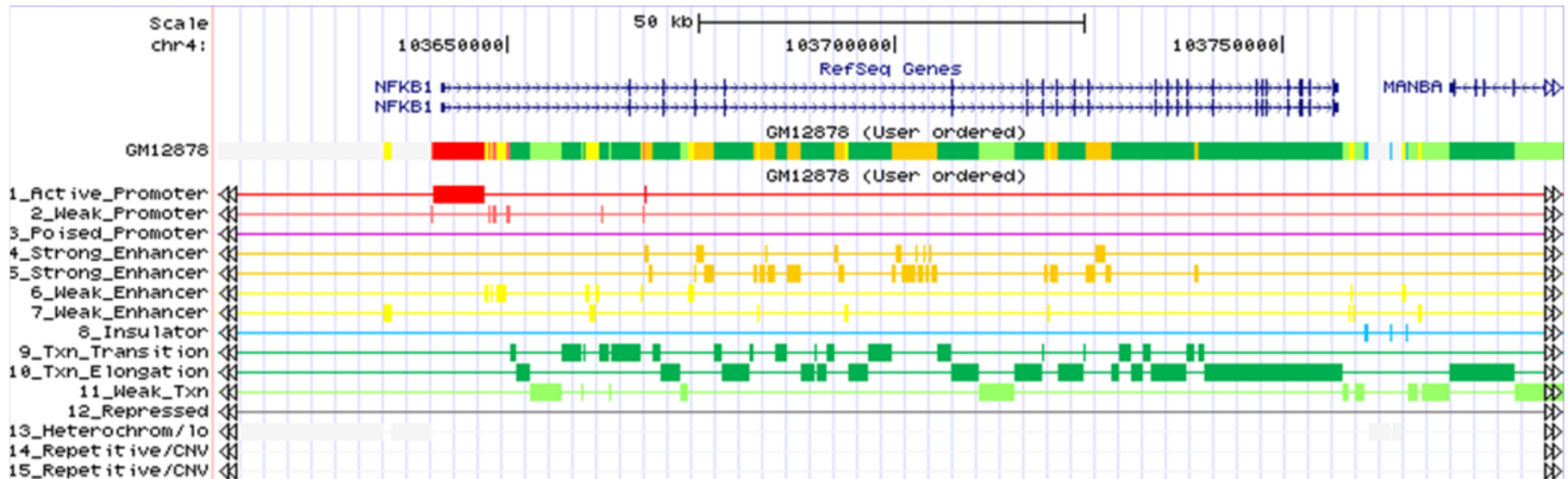# Multi-omics analysis

BMS  Bio-Medical Science Co., Ltd.

# THANK YOU.

**BMS**